# Fashion Retrieval with Contrastive Learning

Lê Hoàng Bùi
*Faculty of Information Technology*
*University of Science, VNU-HCM*
Ho Chi Minh City, Vietnam
0009-0003-1162-234X

Hoang Nguyen Cao
*Faculty of Information Technology*
*University of Science, VNU-HCM*
Ho Chi Minh City, Vietnam
0009-0003-0459-730X

Hoang-Nam Vo
*Faculty of Information Technology*
*University of Science, VNU-HCM*
Ho Chi Minh City, Vietnam
0009-0000-2835-3711

*Abstract—*

*Index Terms—*

## I. INTRODUCTION

## II. RELATED WORK

Composed image retrieval has been an active area of research, with various approaches proposed to address the challenges. Early methods relied on traditional image processing techniques, while recent advancements leverage deep learning. Notable works include the use of convolutional neural networks (CNNs) for feature extraction and the application of contrastive learning for self-supervised learning.

Several approaches have been proposed to address composed image retrieval. Baldrati et al. [1] introduced vision language pre-trained encoders as the backbone. Similarly, Vo et al. [2] developed a cross-modal fusion of text and images input using a gates mechanism. Our work differs from these approaches in using a unified transformer-based architecture for fusing or composing vision-language input, as proposed by [3]. Additionally, progressive learning for image retrieval with hybrid-modality queries was introduced by [4], using self-supervised adaptive weighting mechanisms to determine whether the text or image is more important in a query.

Additional works include the use of multi-modal transformers for composed query image retrieval [3] and the application of sentence-level prompts to benefit composed image retrieval [5].

Data Generation for CIR: InstructPix2Pix [6] first uses GPT-3 to generate modified text for captions and then utilizes a diffusion model to generate images for these texts. CompDiff [?] constructs triplets for CIR datasets by automatically generating modified texts and corresponding images using large language models and diffusion models.

Negative Sampling in Contrastive Learning: SimCLR [?] introduced self-supervised contrastive methods to make similar features "closer" and dissimilar features "further". Unsupervised Feature Learning via Non-Parametric Instance Discrimination [?] focuses on instance-level discrimination, where the features for each instance are stored in a discrete memory bank, rather than weights in a network.

## III. METHODOLOGY

Our framework builds upon the baseline model by introducing enhancements to the query encoder and target image encoder. The query encoder fuses the reference image and modified text to generate a query representation. The target image encoder generates representations for candidate images. We use contrastive learning to optimize the model, focusing on selecting the correct and sufficient amount of positive and negative examples. The objective is to maximize the similarity between the query representation and the target image representation while minimizing the similarity with negative samples.

### A. Problem Formulation

The problem of composed image retrieval involves retrieving a target image based on a reference image and modified text. The goal is to generate a query representation that effectively combines the reference image and modified text to accurately retrieve the target image from a candidate set.

### B. Baseline

Preliminary: Suppose a CIR dataset consists of $N$ annotated triplets consisting of reference images ($r$), modified text ($t$), and target image ($t$). A candidate image set consists of all reference images and target images. The task is to use reference image $r(i)$ and modified text $t(i)$ to compose a query $q(i)$. Then this $q(i)$ is used to retrieve target images from the candidate set.

Paradigm of Composed Image Retrieval: Multiple annotated triplets are combined into a mini-batch, and the reference images and modified texts in the same batch are then encoded using a query encoder to obtain query representations. The target images are encoded using an image encoder to obtain target image representations. The cosine similarity is then adopted to calculate the similarity between the query and target image representations.

The use of contrastive learning: Treat the annotated examples as positive examples and treat the examples obtained by replacing the target image in the positive examples with another image in the mini-batch as the negatives. This approach pulls the query and target representations in positive examples closer while pushing the query and target representations in negative samples further. This yielded good results; however, the lack of positive examples and negative examples severely limits the full performance of contrastive learning, leading to optimization.

Scaling positive examples: Using multi-modal Large Language Model (MLLM) to construct the triplets for CIR. Caption Generation: Design a prompt template to guide the MLLM to obtain a brief caption for each image under constrained conditions, where type and $k$ are two dataset-specific parameters to simulate the type and length of modified text in the real dataset. Image Pair Matching: Match two image-text pairs to generate a quadruplet. Introducing a uni-modal image encoder to get the representation of every image and calculate the pairwise similarity between two different images. Then we can rank the similarities related to the reference image in descending order. Only one image whose similarity rank is between $[c_0, c_1]$ $(c_0 < c_1)$ will be chosen as the target image, where $c_0$ and $c_1$ are two hyperparameters. This forms quadruplets. Modified text generation: From the quadruplets, design a prompt template that generates modified text by either one of these:

- $P_{\text{temp0}} : \{C_t^i\}$ instead of $\{C^i\}$
- $P_{\text{temp1}} :$ Unlike $\{C^i\}$, I want $\{C_t^i\}$
- $P_{\text{temp2}} : \{C_t^i\}$

Positive Example Construction: Combine image pairs from the second step with the modified text obtained in the third step to get new $M$ triplets.

Scaling Negative Examples: Construct negative examples by replacing any element in the triplet. There are four methods to construct negative examples, and the authors examined each method to conclude that replacing the target image will obtain the best negative examples.

Two-stage Finetuning: In the first stage, fine-tune both the query encoder and the target encoder with in-batch negative sampling. In the second stage, freeze the target encoder and only fine-tune the query encoder.

### C. Potential Approaches

Potential approaches to improve composed image retrieval include exploring advanced contrastive learning techniques, incorporating additional data augmentation strategies, and leveraging multi-modal large language models to enhance the quality of modified text descriptions.

## IV. EXPERIMENTS

We conducted extensive experiments to evaluate the performance of our proposed framework. The experiments were carried out on several benchmark datasets. We compared our method with state-of-the-art approaches and measured the retrieval accuracy using metrics such as precision, recall, and F1-score. The results show that our framework outperforms existing methods in terms of retrieval accuracy and robustness.

### A. Datasets

We evaluated our method on two commonly used fashion-domain datasets: *FashionIQ* and *CIRR*.

### B. Evaluation Metrics

We used *Recall@K (R@K)* as the primary evaluation metric, which is the proportion of queries for which the retrieved top $K$ images include the correct target image. *Recallsubset@K (Rsubset@K)* is similar to R@K but the model only retrieves inside the semantically similar group of the reference image.

### C. Ongoing Experiments

We are currently conducting experiments to evaluate the impact of different contrastive learning techniques on the retrieval accuracy. Preliminary results indicate that incorporating additional negative samples significantly improves performance.

## V. CONCLUSION

### REFERENCES

[1] A. Baldrati, M. Bertini, T. Uricchio, and A. Del Bimbo, "Conditioned and composed image retrieval combining and partially fine-tuning clip-based features," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, Jun. 2022, pp. 4959–4968.
[2] N. Vo, L. Jiang, C. Sun, K. Murphy, L.-J. Li, L. Fei-Fei, and J. Hays, "Composing text and image for image retrieval - an empirical odyssey," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2019.
[3] Y. Xu, Y. Bin, J. Wei, Y. Yang, G. Wang, and H. T. Shen, "Multi-modal transformer with global-local alignment for composed query image retrieval," *Trans. Multi.*, vol. 25, no. 1, pp. 8346–8357, Jan. 2023. [Online]. Available: https://doi.org/10.1109/TMM.2023.3235495
[4] Y. Zhao, Y. Song, and Q. Jin, "Progressive learning for image retrieval with hybrid-modality queries," 2022. [Online]. Available: https://arxiv.org/abs/2204.11212
[5] Y. Bai, X. Xu, Y. Liu, S. Khan, F. Khan, W. Zuo, R. S. M. Goh, and C.-M. Feng, "Sentence-level prompts benefit composed image retrieval," 2023. [Online]. Available: https://arxiv.org/abs/2310.05473
[6] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," 2023. [Online]. Available: https://arxiv.org/abs/2211.09800